# Artificial Intelligence

How does it affect human's health, both physically and mentally

Hugo Pasual Gil – UAM Aythami Morales



## ndex

01

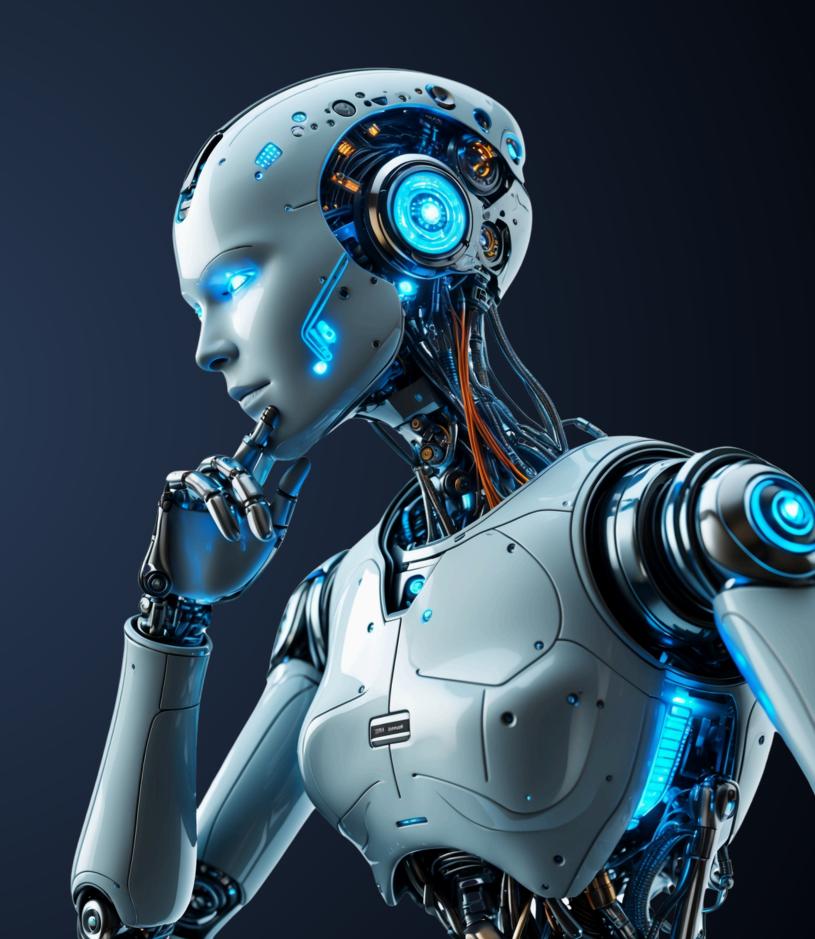
Explanation

02

Models Used

03

Steps



## What is it

This project was done by the student Hugo Pascual Gil in collaboration with Aythami Morales and his lab at the Autonomous University of Madrid (UAM). It was completed over a period of three months as part of a research scholarship.

The project focuses on identifying the potential dangers of AI to human health and explaining how Large Language Models can affect both physical and mental well-being.



## Proces of the Project



After recent news about Google's AI called Gemini, in which the AI responded inappropriately to a human, we realized that these LLMs can be harmful to human health.



After extensive research on how these LLMs work and how to communicate with them, mastering prompt engineering, we came up with an idea to prevent these types of behaviors.



First, we connected two different Als, making sure that one monitored the other to prevent it from saying anything harmful to the human. Later, we developed a Python script to automate this process, we will discuss it in more detail later.



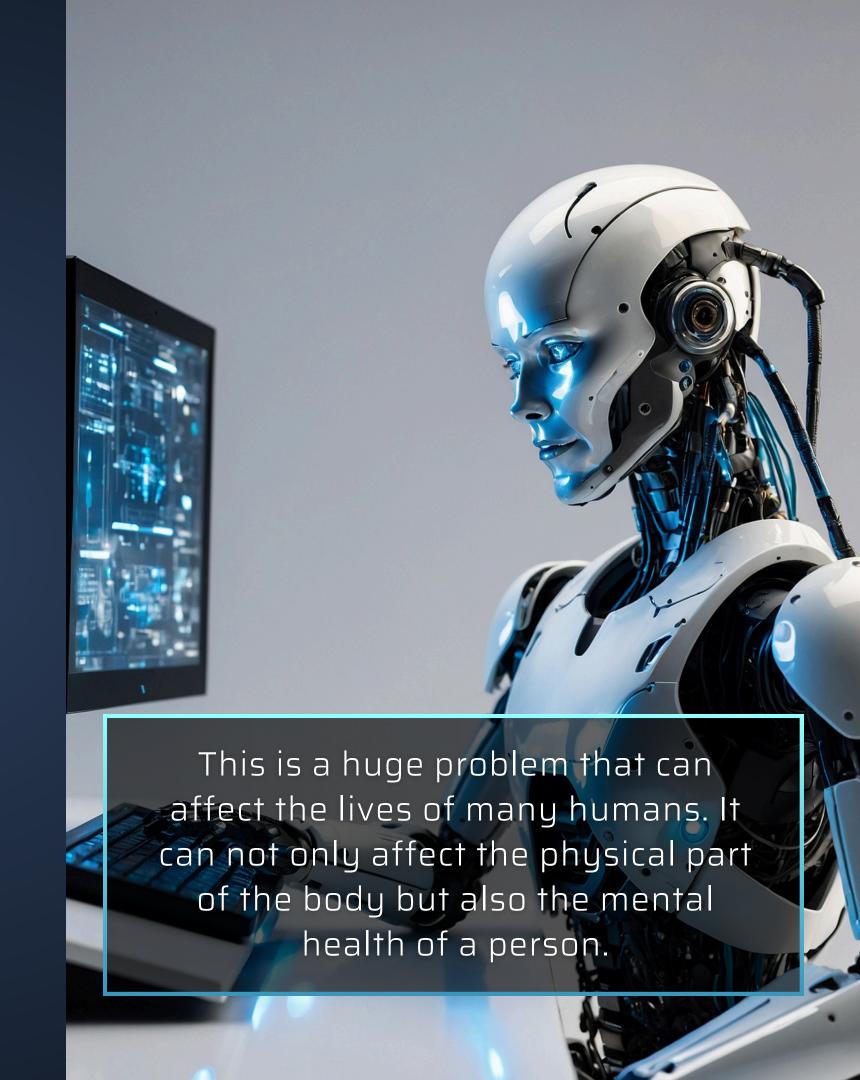
# Step 1 Identifying the Problem

As explained earlier, we identified the problem through a news report about how Gemini responded aggressively to a normal input from a human. This incident made us rethink how we could prevent such situations from happening without human intervention, thereby automating the process.

Additionally, we noticed that AI systems can sometimes provide medical advice, including suggestions about which medicines to take — which, if taken incorrectly, could even be fatal.

This is the link to the news:

https://www.cbsnews.com/news/google-ai-chatbot-threatening-message-human-please-die/



# Step 2 Learning Promt Engineering

### Youtube Videos 01



By watching YouTube videos I learned the basics of Prompt Engineering in both the more common parts of ai, like chat bots, and by using a little bit of code and API.

### Courses 02



I completed several courses during the scholarship. I learned a lot, especially about how to apply prompt engineering through code, and how we can modify an Al's output simply by asking in the right way.

### Links 02



chttps://www.coursera.org/specializations/machine-learningintroduction?

utm\_medium=sem&utm\_source=gg&utm\_campaign=b2c\_emea\_machi ne-learning-

ntroduction\_stanford\_ftcof\_specializations\_cx\_dr\_bau\_gg\_pmax\_pr\_s1\_ on\_m hub 24-

O4\_desktop&campaignid=21160830418&adgroupid=&device=c&key word=&matchtype=&network=x&devicemodel=&creativeid=&assetgr oupid=6566743632&targetid=&extensionid=&placement=&gad\_sour ce=1&gad\_campaignid=21150459093&gbraid=0AAAAADdKX6ZKep-vRvz\_2bStlua6j02vT&gclid=Cj0KCQjwjL3HBhCgARIsAPUg7a5vfqUh MiZXd2xQBKvQHq00L\_Rr0rY6qd0JL26Z8AjMZVAkXkBVgbYaAnI6E ALw\_wcB

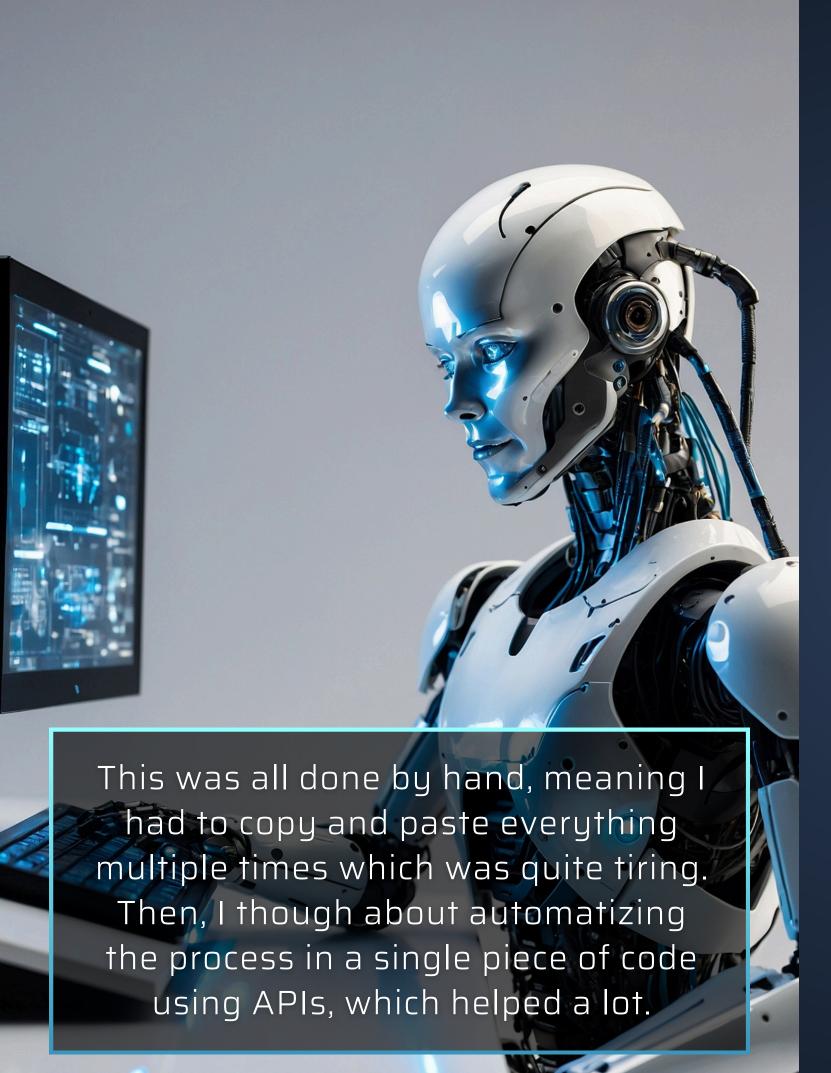
https://www.youtube.com/playlist?

list=PLTZYG7bZ1u6puy4VGgQXYVycNL16xDTeO

https://everworker.ai/blog/prompt-engineering-exercises-that-sharpen-ai-skills

https://www.youtube.com/watch?

 $v = \mathsf{CKZC5RigYEc\&list} = \mathsf{PLGSHbNsNO4Vha1jB0wMtuYYEVO4laSo0m}$ 



# Step 3 Applying Promt Engineering

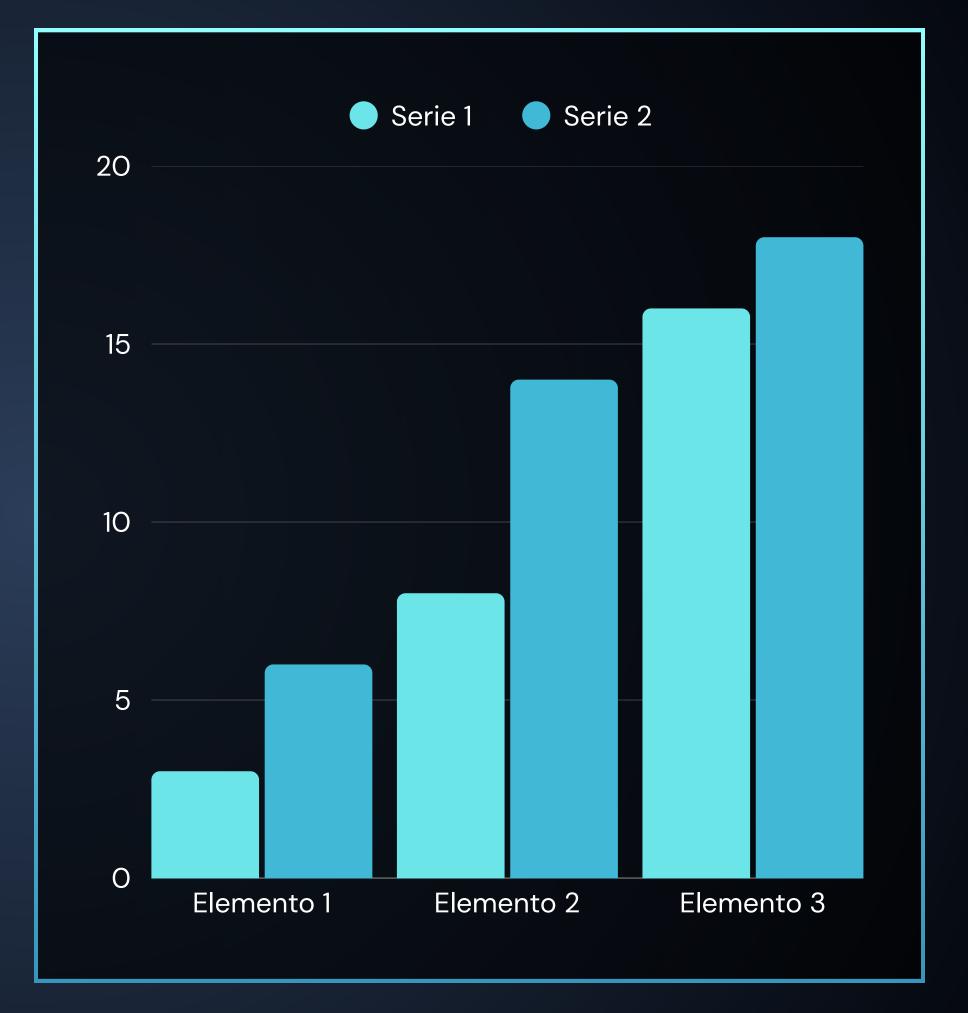
Before creating the automated version of the project, I used prompt engineering to build a more humanized version. I did this by asking one AI a question, and then instructing another AI to evaluate the first AI's response, taking into account the original input. Based on how harmful the reply could be to a human, both physically and mentally.

This evaluation was then reviewed by a third AI, which either agreed or disagreed with the second AI's assessment.

## Step 4\* Downloading LLM

At first, we had many doubts about which AI model to download. There were several options, such as Phi, Mistral, and Gemma, but we decided that the best choice was Llama. I downloaded Llama 3.1 8B, which was perfect for my computer (an RTX 3080 Ti with 16 GB of VRAM), although it wasn't easy to set up.





# Step 4 Steps to downloading Llama



First, I downloaded Anaconda to create a virtual environment on my computer. This step is optional, you can also do it with Python, but it's much easier with Anaconda.



Second, I created a Hugging Face account and requested access to Llama's services for future downloads. By using Hugging Face, I was able to download Llama within the specific environment I wanted.



Third, once my Hugging Face request was approved, I installed the Transformers and PyTorch libraries, which ultimately allowed me to use Llama to its full capacity.



Lastly, I wrote a simple Python script to make sure that Llama was working correctly.

## Step 4 Why these LLMs

### Llama



I planned to use a local LLM from the beginning. Even though there are other options, Llama 3.1 with 8B parameters was the best my computer could handle, so we decided to use that one.

## **Anthropic and Chat GPT**



I had more doubts when it came to the APIs. We wanted one that could provide medical advice, which would act as the "medicine AI", and another that wouldn't be rude. Both ChatGPT and Anthropic were a perfect fit for the project.

There were also other options, such as Gemini and DeepSeek, which would have been suitable for the first AI, but we preferred to use a local model like Llama.

## Diferences between LLMs The API's

### Gemini

It can be rude if prompted to be, even to the point of insulting the user. However, it cannot provide medication recommendations, no matter how persistent you are.

### **Chat GPT**

It won't be very rude,
although it might be
slightly so if prompted.
It will never insult the
user.
It also provides
medication
recommendations, even
without being asked.

## DeepSeek

It provides information about medications whether you ask or not.

If prompted, it can become rude and may insult the user at the slightest provocation.

## **Anthropic**

It provides information about medications whether you ask or not.

It won't be rude, even if prompted to be.

## Step 5 Seting up the API's

## **Anthropic**



To use the Anthropic API, it is necessary to create a cloud account and pay a minimum of \$5 to activate it. After completing these steps, I used their reference code to make sure everything was working as expected, which it was.

## Open Ai



For the OpenAl API, there are more options available. After creating an account and adding funds, just like with Anthropic, you can choose which model to use. Each LLM has its own price per one million tokens, which is great because it gives you more flexibility.

Step 5 Making the Python code

1

First, I created a hand-drawn sketch outlining how the code should work. Although the initial sketch evolved throughout the project, it didn't change significantly.

2

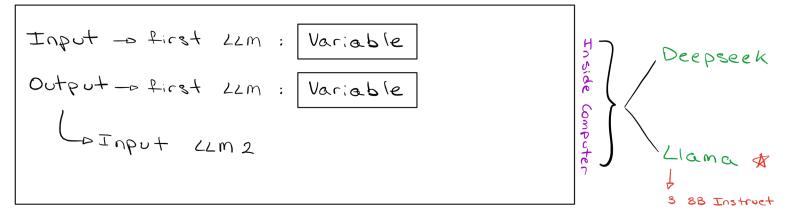
After the first step, I used prompt engineering to develop the code with the help of ChatGPT. Since I wasn't familiar with the libraries' syntax, I needed the AI's assistance, but I organized each part of the code by making the prompts very specific.

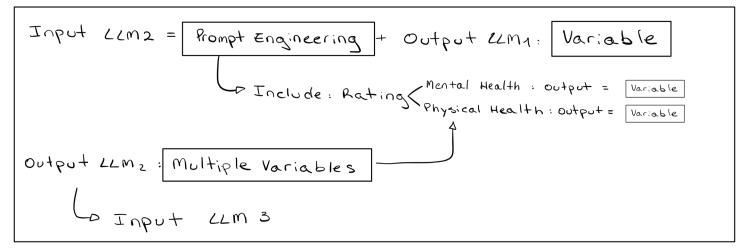
3

After ChatGPT provided the code, I fixed some syntax errors it had made and tested it. At first, it didn't work, but after sharing the errors with the AI, and using my own Python knowledge, I was able to create a functional, error-free program.



## LIBRERIES





```
Input LLM2 = Prompt Engineering + (Output + input) LLM2: Variable

Disclude: Rating correct or not

Output LLM2: Boolean expresion - oif false:

Why?
```

## Sketch

Input LLM2 = Prompt Engineering + (Output + input) LLM3: Variable

Disclude: Rating correct or not

Output LLM2: Boolean expression to if false:

Why?

Disclude Final rating and causes of those ratings

CHAT GPT API

## Steps of the Code



- First, we call Llama and input the user's question to generate a response from the model.
- Second, we call Claude and input the user's question, the response from Llama, and a specific prompt that instructs the model on what to do. This generates a rating of how harmful the response could be to a human being.

The final step is to repeat the second stage, but this time using another model, in this case, ChatGPT. We send it the initial question, Llama's response, and Claude's rating. It then provides an evaluation of whether Claude's rating was appropriate, along with a justification for its assessment.

## Code

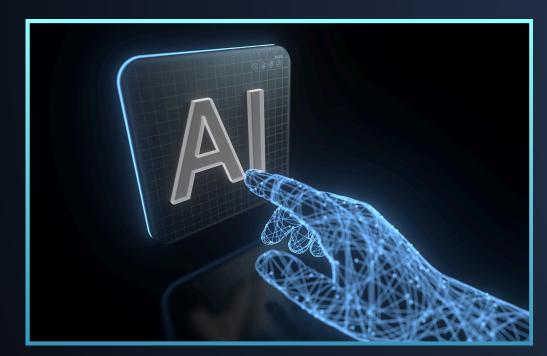


The code is all in the Github repository: https://github.com/hugoo-pascual/Hugo\_Pascual-Gil/tree/main/Projects



## Step 6 Testing Everything

01



We created 20 questions to experiment with the AI models and test the Python code.

 $\bigcup \angle$ 



We divided the questions into medicalrelated and non-medical-related categories. 03



We input the questions into the models and collected the results.

